# DATABASE MANAGEMENT

## CASSIOUS KABWE

### February 14, 2024

1. Explain the Need for Sorting:

   - Summarize the reasons for sorting data as outlined in the provided information.
     - Data Retrieval: When users request data in a specific order (e.g., finding students sorted by increasing GPA), sorting ensures efficient retrieval. Sorted data allows for faster search and retrieval operations.
     - Eliminating Duplicate Copies: Sorting helps identify and eliminate duplicate records. This is essential for maintaining data integrity and ensuring accurate results in queries.
     - Sort-Merge Join Algorithm: In query processing, the sort-merge join algorithm involves sorting data from multiple tables before merging them. Efficient sorting improves join performance.
   - Provide examples of scenarios where sorting is necessary.
     - Search and Retrieval: When users query a database for specific records (e.g., finding all customers who made purchases in the last month), sorting ensures efficient retrieval. Sorted data allows for faster search operations.
     - Join Operations: When combining data from multiple tables using join operations (e.g., INNER JOIN, LEFT JOIN), sorting based on join keys (common columns) is essential. It ensures that matching records align correctly.
     - Aggregations and Grouping: Aggregating data (e.g., calculating total sales, average scores) often involves grouping records by specific attributes (e.g., department, category). Sorting ensures that grouped data is organized logically.

2. Challenges of Sorting Large Datasets:

   - Discuss the challenges associated with sorting large datasets when the available RAM is limited.
     - Memory Constraints: When the dataset exceeds the available RAM, we cannot load the entire dataset into memory for sorting. This limitation affects the choice of sorting algorithms and strategies.

- Slower Performance: Sorting algorithms that rely heavily on random access memory (e.g., quicksort) may perform poorly when constrained by limited RAM. Frequent disk I/O operations slow down the sorting process.
- Choosing the Right Algorithm: Selecting an appropriate sorting algorithm becomes critical. Some algorithms (e.g., merge sort, heap sort) work well with limited RAM, while others (e.g., quicksort) may struggle.

- Explain why virtual memory is not a viable solution for sorting large datasets.

  - Fragmentation: Virtual memory can cause fragmentation—data scattered across non-contiguous memory locations. Fragmented data complicates sorting algorithms and increases access time.
  - Limited Address Space: Virtual memory has a finite address space. Sorting extremely large datasets may exceed this limit, causing addressing issues.
  - Disk I/O Bottleneck: When sorting large datasets, virtual memory heavily relies on disk reads and writes. Disk access is much slower than RAM access, leading to poor performance.

3. Two Themes of Out-of-Core Algorithms:

- Describe the two themes of out-of-core algorithms for handling large datasets.

  - Incremental Processing: In this approach, the dataset is processed incrementally, one chunk at a time. The algorithm reads a portion of the data, performs computations, and updates its model or results.
  - Block-Based Processing: In this approach, the dataset is divided into smaller blocks or chunks. Each block is processed independently, and the results are combined. The data is split into fixed-size blocks (e.g., chunks of rows or columns). while each block is processed concurrently using parallel computing techniques.

4. Two-Way External Merge Sort:

- Explain the process of two-way external merge sort, including its two phases: conquer a batch and merge via streaming.

  - Conquer a Batch: In this phase, the input data is divided into smaller blocks or runs (chunks) that can fit into memory. Each run is sorted in memory using an internal sorting algorithm (e.g., quicksort, insertion sort). Its key feature is to read a chunk of data (run) from the input file, Sort the run in memory, write the sorted run to a temporary file on disk and repeat these steps until all data is processed.

– Merge via Streaming: After creating sorted runs, the goal is to merge them into a single sorted output file. This is done using a two-way merge approach. Its Key steps is to maintain two input buffers (one for each sorted run) and an output buffer, read the first element from each run into the input buffers, compare the elements in the input buffers and select the smaller one and write the selected element to the output buffer.

- Provide a diagram illustrating the process of two-way external merge sort.

5. General External Merge Sort:

- Describe the general external merge sort algorithm, focusing on how it utilizes multiple buffer pages.
  – External merge sort is essential when the entire dataset cannot fit into memory. The goal is to sort a large file that exceeds available memory capacity.We break down the data into smaller chunks (sorted runs) and merge them in subsequent passes. For example, if we have a 108-page file and 5 buffer pages, we create 22 sorted runs of 5 pages each (with the last run having 3 pages).

- Explain the process of sorting a file with N pages using B buffer pages.
  – Utilize B buffer pages to read the input file. Divide the input into (N/B) sorted runs, each containing B pages. Write these sorted runs to disk. The general idea is to divide and conquer: sort smaller chunks of data (runs) and then merge them to create larger sorted runs. The process continues until a single sorted file is obtained.

6. Cost of External Merge Sort:

- Calculate the number of passes and total I/Os required for external merge sort given different values of N and B.
  – In a one-pass system, we perform a single merge of all the chunks while In a two-pass system, we break the data into A groups of B chunks and then perform a final merge. We observe that a two-pass system only reaches the same I/O cost as one pass when the chunk size ( B ) approaches infinity. Even then, the additional final merge still gives us the same I/O cost as the one-pass approach.

- Provide a table showing the number of passes required for various values of N and B.

| buffer-size | number-of-passes |
|:---:|:---:|
| 1 | 1 |
| 2 | 2 |
| 3 | 2 |
| 4 | 2 |

7. Memory Requirement for External Sorting:

- Discuss the memory requirement for sorting a table in exactly two passes using external sorting.

  – Two-Pass External Sorting: In this approach, we divide the data into smaller chunks (runs) that fit into memory (buffers). The goal is to minimize the number of I/O operations while efficiently sorting the entire dataset. Pass one Divide the input data into A groups of B chunks (where each chunk fits into memory), Sort each group of chunks in memory and Write these sorted runs to temporary files on disk while pass two Perform a merge of the sorted runs from the temporary files, use B-1 buffers for merging and write the merged output to the final sorted file.

- Explain how the size of the table affects the memory requirement.

  – External sorting is used to sort large datasets that cannot fit entirely into memory (RAM). It involves dividing the data into smaller chunks (runs) that fit in memory and then merging these runs to create a sorted output. The size of the table (number of pages or records) directly affects the memory needed during external sorting and If the table is larger (more pages or records), we need more memory to hold the chunks during sorting.

8. Alternative: Hashing:

- Describe the concept of hashing as an alternative to sorting.

  – Hashing involves mapping data (such as keys or values) to fixed-size representations (hash codes or hash values). A hash function takes an input (data) and produces a unique hash code. The hash code is typically an integer or a fixed-length string. The advantages of hashing are Hashing allows direct access to data based on its hash code, Hashing distributes data across buckets, useful for distributed systems.

- Explain how hashing can be used to eliminate duplicates and form groups.

  – Hash Set: Stores unique elements (no duplicates). Uses a hash function to map each element to a unique hash code. If an element with the same hash code already exists, it is considered a duplicate and not added.

9. Divide and Conquer in Hashing:

- Explain the two phases of hashing: divide (streaming partition) and conquer (rehash).

  – Divide (Streaming Partition): In the divide phase, we focus on breaking down the data into manageable partitions based on

hash values. Assume we have B buffer frames available (where each frame can hold a page of data). We start by streaming data from disk into our input buffer (one of the B buffers). The goal is to create partitions: A partition is a set of pages where all values on those pages hash to the same hash value (determined by a specific hash function). Each partition corresponds to one of the B-1 output buffers.

– Conquer (Rehash): In the conquer phase, we construct the final hash tables from the partitions. For partitions that can fit in memory (i.e., their size is less than or equal to B pages), we proceed with building the hash tables. The process involves: Reading the pages of each partition into memory. Constructing an in-memory hash table for each partition. Ensuring that all occurrences of a value are in the same partition (due to the hash function properties). Once all partitions are processed, we concatenate the hash tables to form the final result.

• Provide examples illustrating the process of divide and conquer in hashing.

– Merge Sort: Merge Sort is a classic example of divide and conquer first its to divide which means to split the array into two halves then conque which is to recursively sort each half. and lastly Combine which involves merging the sorted halves to create the final sorted array.

– Anagrams Grouping: Suppose we have a list of words. Divide can happen if we hash each word based on its letter frequencies (e.g., count of each letter) and conquer can happen in the case that for each partition, sort the anagrams alphabetically.

– External hashing requires less memory (as much as possible) while External sorting uses more memory due to sorting operations.

10. Cost of External Hashing:

• Compare the cost of external hashing with external sorting for the one-pass case.

– Memory Requirement: External hashing requires less memory (as much as possible) while external sorting uses more memory due to sorting operations.

– I/Os: Both have similar I/O costs (read and write passes) while external hashing may have a slight advantage due to fewer passes.

• Discuss the factors influencing the cost of external hashing.

– Buffer Size (B): The number of buffer pages available for reading and writing during hashing significantly impacts the cost.

– Hash Function Efficiency: The quality of the hash function used affects the distribution of hash codes.

- Memory Constraints: Too little memory leads to frequent disk access, while too much memory may impact other system processes.

11. Parallelization in Sorting and Hashing:

- Describe how sorting and hashing algorithms can be parallelized across multiple machines.
    - Parallelizing Sorting Algorithms: Many sorting algorithms (e.g., merge sort, quicksort) follow a divide-and-conquer approach to divide the dataset into smaller subproblems.
    - Parallelizing Hashing Algorithms: Some hash functions (e.g., SHA-256) can be computed in parallel and to divide the input data into chunks and hash them concurrently.
- Explain the phases involved in parallelization for both sorting and hashing.
    - Divide and Conquer: Many sorting algorithms follow a divide-and-conquer approach.
    - Parallel Merge Sort: Divide the array into chunks then sort each chunk in parallel using separate machines and finally Merge the sorted chunks to create the final sorted array.
    - Parallel Quick Sort: Choose a pivot element, partition the data into smaller segments, sort each segment in parallel and combine the sorted segments.
    - Parallel Radix Sort: Divide the data into buckets based on digits, sort each bucket in parallel and combine the buckets to get the sorted result.

12. Summary:

- Summarize the key points discussed in the reading material, including the sort/hash duality, the benefits of sorting and hashing, and the importance of memory hierarchy.
    - Sort/Hash Duality: Sorting arranges data in a specific order (e.g., ascending or descending) while hashing maps data to fixed-size representations (hash codes)
    - Importance of Memory Hierarchy:
        * Cache: Utilize cache memory for faster data access.
        * RAM: Optimize memory usage for sorting and hashing.
        * Disk: Minimize I/O operations by managing data efficiently.
        * External Storage: Use SSDs or other external storage for large datasets.
    - Benefits of Sorting:
        * Ordering: Sorting organizes data for easier search, retrieval, and analysis.

* Search Efficiency: Sorted data allows binary search and other efficient search algorithms.
* Data Presentation: Sorted data is visually appealing and easier to understand.
- Benefits of Hashing:
  * Fast Access: Hashing provides direct access to data based on hash codes.
  * Distributed Data Storage: Hashing distributes data across buckets (useful in distributed systems).
  * Data Integrity: Cryptographic hash functions ensure data integrity.